

This manuscript is a post-print copy of the following article

Title: **P-N-RMiner: A generic framework for mining interesting structured relational patterns**

Authors: Jefrey Lijffijt ; Eirini Spyropoulou ; Bo Kang ; Tijl De Bie

The final publication is available at IEEE via <https://doi.org/10.1109/DSAA.2015.7344837>

© 2015 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

# P-N-RMiner: A Generic Framework for Mining Interesting Structured Relational Patterns

Jefrey Lijffijt, Eirini Spyropoulou, Bo Kang, Tijl De Bie  
Intelligent Systems Lab  
University of Bristol, UK  
Email: jefrey.lijffijt@bristol.ac.uk

**Abstract**—Local pattern mining methods are fragmented along two dimensions: the *pattern syntax*, and the *data types* on which they are applicable. Pattern syntaxes considered in the literature include subgroups,  $n$ -sets, itemsets, and many more; common data types include binary, categorical, and real-valued. Recent research on pattern mining in relational databases has shown how the aforementioned pattern syntaxes can be unified in a single framework. However, a unified understanding of how to deal with various data types is lacking, certainly for more complexly structured types such as time of day (which is circular), geographical location, terms from a taxonomy, etc.

In this paper, we introduce a generic approach for mining interesting local patterns in (relational) data involving such structured data types as attributes. Importantly, we show how this can be done in a generic manner, by modelling the structure within a set of attribute values as a partial order. We then derive a measure of subjective interestingness of such patterns using Information Theory, and propose an algorithm for effectively enumerating all patterns of this syntax. Through empirical evaluation, we found that (a) the new interestingness derivation is relevant and cannot be approximated using existing tools, (b) the new tool, P-N-RMiner, finds patterns that are substantially more informative, and (c) the new enumeration algorithm is considerably faster.

## I. INTRODUCTION

**Motivation.** Local pattern mining, including itemset mining and variants [1],  $n$ -set mining [2], subgroup discovery [3], and multi-relational pattern mining [4], [5], has traditionally been rooted in categorical (or even binary) data. Some of these local pattern mining approaches have been extended in various ways to include ordinal, interval (e.g. real-valued), or otherwise structured data. For example, extensions of itemset mining for real-valued data has led to approaches akin to biclustering, and subgroup discovery methods that allow discovery of rules based on attribute-value inequalities are available.

This work is fragmented and often ad hoc, in the sense that other kinds of structure (e.g., taxonomy terms, time-of-day intervals on a circular 24-hour clock, geographical regions on the globe, etc.) may not be approachable in the same way and may necessitate fundamentally different approaches. The purpose of this paper is to provide an elegant and encompassing framework to deal with attributes of any of the structured types listed above and more, and this in a relation setting (i.e., applicable to data as it resides in relational databases). To illustrate the breadth and nature of the contributions, we provide two motivating examples.

*Example 1:* Consider a dataset of Foursquare<sup>1</sup> check-in times of a number of users. Such a dataset has the potential of elucidating particular lifestyle patterns shared by a number of Foursquare users. To formalise and then find such patterns, it is tempting to specify a time resolution and discretise the data. However, it is unclear which discretisation level to use, and whether to take it uniform throughout the day. In fact, the optimal discretisation could vary for various lifestyle patterns.

An alternative approach could be to take the mean and possibly higher-order statistics of the check-in times for each user, and find patterns in this summary description. This approach would suffer from two problems: first, computing averages of circular quantities is ambiguous (e.g. is the mean of 6am and 6pm midnight or noon?), and second, it ignores much of the information in the data.

The method developed in this paper, when applied to this data, deems as most interesting a pattern that reveals that 1.6% of all users check in frequently in the 6am-7am interval and again in the 10.10am-10.50am interval. Here, the interval sizes are tuned automatically to maximise interestingness.

While the first example illustrates how the contribution in this paper advances the state-of-the-art even for a single relation (between users and check-in times), the second example shows the full power on data in a relational database.

*Example 2:* Consider a relational database involving users, who have rated books (with an integer from 1 to 5), which are tagged with a number of genres organised in a taxonomy. Applied to this dataset, the method proposed in this paper identifies interesting patterns in the form of sets of books that have been rated by the same set of users in a similar way (say, in the 3-5 interval), which may all belong to a particular set of genres (say, fantasy and action).

This second example illustrates the ability of the proposed method to identify patterns that span several types of entities (users, ratings, books, genres), including structured ones with e.g. ordinal values or values that are organised in a taxonomy.

**Contributions.** The work in this paper is most easily explained as an extension of the N-RMiner algorithm for mining local patterns in relational databases [6], towards structured entity types. However, given the generality of the N-RMiner pattern syntax, this immediately results in a method that includes itemset mining,  $n$ -set mining, and subgroup discovery for structured data types as special cases. To do this, we overcome the following challenges.

---

<sup>1</sup><https://foursquare.com/>



Figure 1. Example schema of users and check-in times. Additionally, we know the age and profession of the users. There are three relationship types: there are relationships between (1) users and check-in times, (2) users and ages, and (3) users and professions.

- We formalise the problem and a matching pattern syntax, in a manner as generic as possible (Sec. II). To achieve this, we adopt an abstract formalisation in terms of a partial order over the structured values. For example, with the time-of-day and book ratings, the partial order is over the intervals, where one is ‘smaller’ than another if it is included in it. For taxonomy terms, one taxonomy term is smaller than another if it is a specialisation of it.
- We formalise the interestingness of such patterns. This is a non-trivial contribution over the approach applicable for the N-RMiner pattern syntax (Sec. III).
- We provide an algorithm for efficiently enumerating all such patterns. This is a non-trivial extension of the algorithmic approach used in N-RMiner (Sec. IV).

## II. PROBLEM FORMALISATION

**Notation.** We formalise a *relational database* as a tuple  $\mathcal{D} = (E, t, \mathcal{R}, R, \succeq)$ . Here,  $E$  denotes the set of *entities*, and  $t : E \rightarrow \{1, \dots, k\}$  is a function that gives the *type* of an *entity* (assuming  $k$  types).  $\mathcal{R}$  denotes the set of all *relationship instances* in the database, while  $R \subseteq \{1, \dots, k\} \times \{1, \dots, k\}$  denotes the set of tuples of entity types whose entities may have relationships, according to the schema of the database. The elements of  $R$  will be referred to as the *relationship types*. So far, this is identical to the formalisation in [5].

As an example, consider the schema illustrated in Figure 1. There are four entity types: *User* (1), *Check-in times* (2), *Profession* (3), and *Age* (4). The numbering is arbitrary. The set  $E$  contains all entities of all types. The set of allowed relationships is  $R = \{(1, 2), (1, 3), (1, 4)\}$  and  $\mathcal{R}$  contains all actual instances of such relationships.

In the Check-ins data (Figure 1), *Age*, *Check-in times*, and *Profession* could all be structured attributes; the values of *Age* are numerical, *Check-ins times* are numerical but without full order, and *Profession* has hierarchical structure. One could be interested in finding patterns in such data not only including an exact age such as 32, but also intervals such as [25–35]. The set of all such intervals can be modelled as a partial order. An example of such a partial order is given in Figure 2.

Hence, we consider one additional element in the data model: a partial order  $\succeq$  that represents implication of relationships across entities of the same type. That is,  $e \succeq f$  means that if any entity  $g$  is related to  $f$ , i.e.,  $(f, g) \in \mathcal{R}$ , then  $g$  is also related to  $e$ :

$$\forall e, f, g \in E : e \succeq f \wedge (g, f) \in \mathcal{R} \Rightarrow (g, e) \in \mathcal{R}.$$

Only implications between entities of the same type are allowed:  $e \succeq f \Rightarrow t(e) = t(f)$ . We assume that  $\mathcal{R}$  contains both

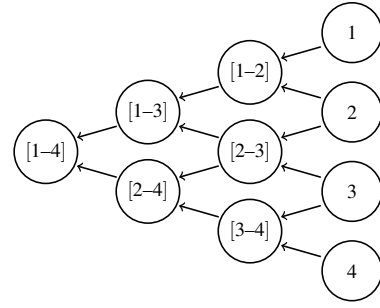


Figure 2. Partial order of all intervals that are supersets of  $\{1\}$ ,  $\{2\}$ ,  $\{3\}$ , and  $\{4\}$ . The partial order corresponds to the superset relation.

the relationship instances between the basic entities present in the database, as well as all relationship instances implied by  $\succeq$ . Where possible, we do not store implied edges explicitly. Hence, we assume  $\mathcal{R}$  contains all relationship instances only for notational convenience.

**Pattern Syntax.** Our aim is to find sets of entities that have surprisingly many relationships. We will refer to a set of entities and the relationship instances among them as a *pattern*. The approach that we take is to first enumerate all possibly interesting patterns, and then rank them according to their interestingness. We define a pattern as potentially interesting if it is *complete*, *connected*, *maximal*, and *proper*.

**Definition 1:** An entity set  $F \subseteq E$  is **complete** iff

$$\forall (t_1, t_2) \in R, \forall e_i, e_j \in F, t(e_i) = t_1, t(e_j) = t_2 : (e_i, e_j) \in \mathcal{R}.$$

**Definition 2:** An entity set  $F \subseteq E$  is **connected** iff

$$\forall e, f \in F, e \neq f : (e, f) \in \mathcal{R} \vee \exists g \in F, \{e, g\} \text{ connected} \wedge \{f, g\} \text{ connected}.$$

**Definition 3:** An entity set  $F \subseteq E$  is **maximal** iff

$$\nexists e \in E : F \cup \{e\} \text{ is complete and connected}.$$

**Definition 4:** An entity set  $F \subseteq E$  is **proper** iff

$$\forall e \in F, f \in E, f \succeq e : f \in F.$$

That is, a pattern  $F$  is *complete* iff all relationship instances between entities in  $F$  that are allowed by the database schema are also present. A pattern  $F$  s.t.  $|F| \geq 2$  is *connected* iff there is a path between any two entities in  $F$  using only entities in  $F$ . Any  $F$  s.t.  $|F| \leq 1$  is *connected*. A pattern  $F$  is *maximal* iff no entity can be added without breaking completeness. *Proper* means that all super entities of any entity in  $F$  are also in  $F$ .

We refer to sets that are complete, connected, and proper as *complete connected proper subsets (CCPSs)*, and to sets that are also maximal as *maximal CCPSs*. In Section IV, we will show that we can enumerate all maximal CCPSs using the so-called fixpoint-enumeration algorithm. The number of such maximal CCPSs may be large though, so in a second step we score and rank them according to an appropriate interestingness measure.

In short, we add a properness constraint and the pattern syntax is otherwise equivalent to [5], [6]. Our implementation and theory also support n-ary relationships, but we do not discuss this further in order to prevent unnecessary complications in the exposition. One could also mine approximate patterns by discarding the completeness constraint. This would lead to an increased computational complexity, but the increase has been shown to be manageable [7]. For simplicity, we do not consider approximate patterns in this paper.

### III. INTERESTINGNESS

**General Approach.** Although we limit the output to maximal CCPs, the number of patterns can be—and often is—exponential in the number of entities. Therefore, it is vital to have a mechanism for identifying the most interesting CCPs. To achieve this, we build upon the framework for interesting pattern mining introduced by De Bie [8], [9]. This framework is based on modelling the user’s prior belief state about the data by means of a Maximum Entropy (MaxEnt) distribution, subject to any stated prior beliefs the user may hold about the data. This distribution is referred to as the *background distribution*. The interestingness of a pattern is then formalised by contrasting the pattern with this background distribution, as the ratio of two quantities:

- the *self-information* of the pattern, defined as minus the logarithm of the probability that the pattern is present under the background distribution, and
- the *description length* of the pattern, which should formalise the amount of effort the user needs to expend to assimilate the pattern.

Given the dependence of this measure on the background distribution, which may in principle differ for different users, this interestingness measure is a *subjective* quantity.

In [5], this framework is used successfully to formalise the interestingness of Complete Connected Subsets (CCSs), without the properness requirement that lies at the core of the contributions in this paper. The properness requirement, however, creates an opportunity as well as a non-trivial challenge. It allows to describe single patterns that capture information that previously could only be presented in a set of patterns. Such patterns reduce the description length. On the other hand, it is harder to compute the self-information of a pattern. We briefly discuss these issues in the next paragraphs, before discussing the latter in greater detail in Sec. III-B.

**Description Length.** The description length of a CCS pattern is formalised as an affine function of the number of entities  $|F|$  in  $F$ . More specifically, with  $|E|$  the total number of entities in the database [5]:

$$\text{DescriptionLength}(F) = |F| \log \left( \frac{1-p}{p} \right) + |E| \log \left( \frac{1}{1-p} \right),$$

where  $p \in (0, 1)$  is a parameter that trades off the cost between describing the presence of an entity in the pattern  $F$  (cost  $\log(p)$ ) and describing its absence (cost  $\log(1-p)$ ).

However, to convey a CCPS pattern to the user, only the *minima* of  $F$  need to be described. Indeed, the presence of the entities that are larger is implied; explicitly describing these would be redundant. Thus, the above expression needs

to be modified by replacing  $|F|$  with the number of *minima* in  $F$ . This leads to a smaller description length than would be required if the partial order  $\succeq$  would be unknown or unaccounted for.

**Information Content.** In the following section, we argue that the background distribution can be fitted in the exact same way as in [5]. However, how to compute the probability that a given pattern is present—and thus its self-information—is not trivial. The difficulty stems from the fact that the presence of relationship instances is now dependent, owing to the partial order relation over the entities. Nonetheless, Sec. III-B describes how the probabilities can still be computed effectively by using the inclusion-exclusion principle.

#### A. The background distribution

In [5], interestingness is formalised under the assumption that users have prior beliefs on the number of entities of a specific type to which a given entity is related. It is argued that this is often a good assumption, and the experiments in the current paper also support that.<sup>2</sup> This assumption leads to a tractable distribution, under which the relationship instances are independent with probabilities that can be found by solving an efficiently solvable convex optimisation problem.

This background distribution factorises over the different relationship types, such that the self-information can be decomposed into a sum of different contributions, each one of which corresponds to the relationship instances for one particular relationship type. That is also the case in the present paper, such that in the rest of this exposition it suffices to imagine just a single relationship type.

What is new is that we implicitly make a further assumption on the user’s knowledge state, namely that the user knows the partial order  $\succeq$ , and hence the fact that if  $e \succeq f \wedge (g, f) \in \mathcal{R}$ , then  $(g, e) \in \mathcal{R}$ . This creates hard-to-handle dependencies between the presence of relationship instances. In practice, data will often only contain relationship instances between minimal entities. In this case, the background distribution can be fitted on the set of minimal entities without worrying about the dependencies, in exactly the same way as done in [5].

In particular, we assume prior beliefs on the number of relationship instances each (minimal) entity is involved in, for every relationship type. The distribution of maximum entropy subject to these prior belief constraints is then used as the background distribution. As shown in [5], this background distribution is a product of Bernoulli distributions, with one factor for each possible relationship instance. In other words: for each possible relationship instance  $(e, f)$ , the distribution gives us a probability  $p_{(e,f)}$  that  $(e, f)$  is present in the data.

This background distribution defines the probabilities  $p_{(e,f)}$  of relationship instances between *minimal* entities  $e$  and  $f$ . Given this, it is possible to compute the probability  $p_{(e,f)}$  of any relationship instance  $(e, f)$ , whether minimal or not, as the probability of presence of any of the relationship instances  $(e', f')$  with  $e \succeq e'$  and  $f \succeq f'$  and  $e'$  and  $f'$  minimal. Indeed, the presence of any such  $(e', f')$  would imply the presence of  $(e, f)$ . How this probability is computed, and how it can be

<sup>2</sup>Of course, exploring other types of prior beliefs is an important line of further work.

used to compute the overall probability of a CCPS pattern given the background distribution, is the subject of Sec. III-B.

More generally, for data that includes relationship instances between non-minimal entities, let us define a partial order  $\succeq_{\mathcal{R}}$  over the relationship instances as follows:  $(e_1, f_1) \succeq_{\mathcal{R}} (e_2, f_2)$  iff  $e_1 \succeq f_1$  and  $e_2 \succeq f_2$ . Then, we suggest fitting the background distribution as before on the minimal relationship instances only. This includes the approach from the previous paragraph as a special case.

This model is imperfect, as the user should be aware of *negative* dependencies between the presence of a relationship instance as a minimal one: if  $(e_2, f_2)$  is a minimal relationship instance, then  $(e_1, f_1)$  with  $(e_1, f_1) \succeq_{\mathcal{R}} (e_2, f_2)$  and  $(e_3, f_3)$  with  $(e_3, f_3) \preceq_{\mathcal{R}} (e_2, f_2)$  cannot be minimal relationship instances. Yet, we argue that in this case, assuming independence is nonetheless still a good approximation.<sup>3</sup>

### B. The self-information of a CCPS

Given a pair of entities  $(e, f)$  such that  $(t(e), t(f)) \in R$  (i.e., they may be related according to the database schema), let us denote the event that  $(e, f) \notin \mathcal{R}$  as  $A_{(e,f)}$  ( $A$  for Absent). The probability of this event under the background distribution can be computed as:<sup>4</sup>

$$P(A_{(e,f)}) = \prod_{(e', f') : (e, f) \succeq_{\mathcal{R}} (e', f')} (1 - p_{(e', f')}).$$

The presence of a CCPS pattern  $F$  corresponds to the event defined by the complement of the union of all events  $A_{(e,f)}$  with  $e, f \in F$  and  $(t(e), t(f)) \in R$ . Hence, the union of all these events corresponds to the event where at least one of the relationship instances is missing. The complement of the union of absence events implies the presence of the pattern. Defining  $\mathcal{T}_F$  as  $\mathcal{T}_F = \{(e, f) | e, f \in F, (t(e), t(f)) \in R\}$ , the set of pairs of entities in  $F$ , this means that the probability of a pattern is given as  $1 - P(\bigcup_{(e,f) \in \mathcal{T}_F} A_{(e,f)})$ . Note that it suffices to consider only the minimal relationship instances from  $\mathcal{T}_F$ , as  $\neg A_{(e,f)}$  implies  $\neg A_{(e', f')}$  for any  $e' \succeq e, f' \succeq f$ .

Directly computing this probability is nontrivial, given the dependencies between the events  $A_{(e,f)}$ . Fortunately, we can use the inclusion-exclusion principle to compute it as follows:

$$P\left(\bigcup_{(e,f) \in \mathcal{T}_F} A_{(e,f)}\right) = \sum_{I \subseteq \mathcal{T}_F} (-1)^{|I|-1} P\left(\bigcap_{(e,f) \in I} A_{(e,f)}\right).$$

Now, the probability of the intersection of events  $A_{(e,f)}$

<sup>3</sup>The intuition is as follows. In practice the probabilities for relationship instances under the background distribution are small. Additionally, for two events with small probabilities  $p$  and  $q$ , the probability of their union is between  $p + q$  (in the case of perfect negative dependence) and  $1 - (1 - p)(1 - q) = p + q - pq$  (in the case of independence), which differs by only  $pq$ , such that assuming independence results in at most a second order error in the probabilities.

<sup>4</sup>As pointed out in Sec. III-A, this expression is exact for databases where relationship instances involve only minimal pairs, and a good approximation in practice in other cases. Note also that only minimal relationship instances have positive probability, and hence non-minimal instances can be ignored.

can be computed straightforwardly as:<sup>5</sup>

$$P\left(\bigcap_{(e,f) \in I} A_{(e,f)}\right) = \prod_{(e', f') : (e, f) \succeq_{\mathcal{R}} (e', f')} 1 - p_{(e', f')}.$$

Hence, we can compute the probability of the presence of a pattern. The self-information is then given as the negative logarithm of this probability:

$$\text{SelfInformation}(F) = -\log\left(1 - P\left(\bigcup_{(e,f) \in \mathcal{T}_F} A_{(e,f)}\right)\right).$$

## IV. ENUMERATION ALGORITHM

Last but not least, we study how to efficiently enumerate all maximal CCPSs. Like previous work on mining interesting patterns in relational data [5], [6], [7], our algorithm is based on the *fixpoint-enumeration algorithm* by Boley et al. [10]. Although that algorithm already exists, it should be noted that it is a meta-algorithm, which does not directly work on the data. The fixpoint-enumeration algorithm takes as input a *set system* and a *closure operator* that together define the problem setting and the output (definitions given below).

We first introduce the fixpoint-enumeration algorithm, after which we introduce notation and formalise our practical problem of enumerating maximal CCPSs as a problem of enumerating all fixpoints in a set system. We prove that the introduced set system is *strongly accessible*, which is required for the fixpoint enumeration to be applicable, and present a suitable closure operator. Finally, we analyse the computational complexity.

**The Enumeration Algorithm.** The fixpoint-enumeration algorithm can efficiently enumerate all fixpoints in a strongly accessible set system  $(E, \mathcal{F})$ , where  $E$  is a set of objects called the *ground set* and  $\mathcal{F} \subseteq \mathcal{P}(E)$  a family of sets. The *fixpoints* are defined by a closure operator  $\sigma$ . The output of the algorithm is valid if and only if the set system satisfies certain criteria [10]. The algorithm is very simple:

- (1) Start with the empty set as the current set:  $F := \{\emptyset\}$ .
- (2) Compute the closure of the current set:  $F := \sigma(F)$ . This closure is one of the fixpoints to return.
- (3) If  $\exists G \supseteq F : G \in \mathcal{F}$ , then pick any element  $f \in G \setminus F : F \cup \{f\} \in \mathcal{F}$  and recurse from (2) to one branch where every set contains  $f$  and one branch where no set contains  $f$ . If there is no such superset then this branch ends.

If the set system  $(E, \mathcal{F})$  is strongly accessible, then all sets in  $\mathcal{F}$  can be found by adding elements one by one while traversing only over sets in  $\mathcal{F}$ . The closure operator defines the fixpoints, which should be interpreted as the subset of sets from  $\mathcal{F}$  that we would like to enumerate.

**Enumerating CCPSs.** The set of all CCPSs forms a set system  $(E, \mathcal{F})$  where the ground set  $E$  is the set of entities and  $\mathcal{F}$  is the set of *valid patterns*, defined as

$$\mathcal{F} = \{F \in \mathcal{P}(E) : F \text{ connected} \wedge F \text{ complete} \wedge F \text{ proper}\}.$$

<sup>5</sup>Note again that only minimal relationship instances  $(e', f')$  need to be considered, since non-minimal relationship instances have zero probability.

The fixpoint-enumeration algorithm can be used to enumerate all closed patterns from this set system, and it is efficient if we can define an appropriate closure operator. Ultimately, we are interested in enumerating the maximal CCPSSs, while  $\mathcal{F}$  contains all CCPSSs.

**Strong Accessibility.** For the fixpoint-enumeration algorithm to be applicable, the set system must be *strongly accessible*. This is the case iff

$$\forall F \in \mathcal{F} \setminus \{\emptyset\} : \exists e \in F : F \setminus \{e\} \in \mathcal{F}, \text{ and} \quad (1)$$

$$\forall F, F' \in \mathcal{F}, F \subset F' : \exists e \in F' \setminus F : F \cup \{e\} \in \mathcal{F}. \quad (2)$$

*Theorem 1:*  $(E, \mathcal{F})$  is strongly accessible.

*Proof:* We prove each of the two properties separately, but first we introduce some notation for convenience. Let  $(F, \succeq)$  denote the set  $F$  partially ordered by  $\succeq$ . We write that an entity  $e \in F$  is minimal in  $(F, \succeq)$  iff  $\nexists f \in F, e \neq f, e \succeq f$ . Likewise an entity  $e \in F$  is maximal in  $(F, \succeq)$  iff  $\nexists f \in F, e \neq f, f \succeq e$ .

The first property states that for every CCPS  $F$ , there should be an entity  $e \in F$  that can be removed such that we obtain another CCPS  $F' = F \setminus \{e\}$ . We prove this by narrowing down candidates by looking in turn at completeness, properness, and finally connectedness:

- (1)  $\forall F \in \mathcal{F} \setminus \{\emptyset\} : \exists e \in F : F \setminus \{e\} \in \mathcal{F}$ , because
  - Removing an entity never violates completeness.
  - Any  $e \in F, e$  minimal in  $(F, \succeq)$  can be removed without breaking properness, and  $\exists e \in F, e$  minimal in  $(F, \succeq)$
  - If  $\exists e, f \in F, e \succeq f, f$  minimal in  $(F, \succeq)$ , then  $F \setminus \{f\} \in \mathcal{F}$ , because  $F \setminus \{f\}$  is complete and proper (see two previous statements) and since  $F$  is connected, for any  $(f, g) \in \mathcal{R}$  also  $(e, g) \in \mathcal{R}$  (since  $e \succeq f$ ), thus  $F \setminus \{f\}$  is also connected.
  - If  $\nexists e, f \in F, e \succeq f, f$  minimal in  $(F, \succeq)$ , then  $\forall e \in F : e$  minimal in  $(F, \succeq)$ . Hence, removal of any entity would not break completeness or properness. Then, we could model the entities of  $F$  as nodes in a graph and the relationship instances between entities in  $F$  as its edges. Since  $F$  is connected, that graph is also connected. Any connected graph has a spanning tree and it is possible to remove any leaf node from that spanning tree without breaking connectedness of the graph.

The second property states that for any pair of CCPS  $F, F' \in \mathcal{F}, F \subset F'$ , there is an entity  $e \in F' \setminus F$  that can be added to  $F$  to lead to another CCPS  $F \cup \{e\} \in \mathcal{F}$ . We prove this property by considering all entity types of entities that are in  $F'$  and not in  $F$ , and then we condition on whether  $F$  is the empty set or whether it already contains some entities.

- (2)  $\forall F, F' \in \mathcal{F}, F \subset F' : \exists e \in F' \setminus F : F \cup \{e\} \in \mathcal{F}$ , because
  - Let  $t(F) = \{t_j | t_j = t(e), e \in F\}$ . For every type  $t_j \in t(F' \setminus F)$ ,  $\exists e \in F' \setminus F : t(e) = t_j, e$  maximal in  $(F' \setminus F, \succeq)$ , since  $F' \setminus F$  is finite.
  - If  $F = \emptyset$ , then for  $\forall e \in F' \setminus F, t(e) = t_j, e$  maximal in  $(F' \setminus F, \succeq) : F \cup \{e\} \in \mathcal{F}$ .
  - If  $F \supset \emptyset$ , then because every entity type has one or more maximal elements and  $F'$  is connected, there is a type adjacent to or present in  $F$  which includes an entity  $e$  maximal in  $(F' \setminus F, \succeq)$  and then  $F \cup \{e\}$  is complete, connected and proper. ■

**The Closure Operator.** Strong accessibility implies that we can efficiently enumerate all fixpoints in  $\mathcal{F}$  in a single traversal over the set system without considering any set twice [10]. A trivial choice for the fixpoints would be all sets in  $\mathcal{F}$ ; in which

case  $\sigma(F) = F, \forall F \in \mathcal{F}$ . However, in the worst case the number of CCPSSs  $|\mathcal{F}|$  is an exponential in  $|E|$ , while there is only one maximal CCPS. Hence, we would like to choose the set of fixpoints such that it includes all maximal CCPSSs and as few other CCPSSs as possible. It is not possible to choose the closure operator such that we enumerate only maximal CCPSSs, because a CCPS may have multiple maximal extensions.

We derive a suitable closure operator from its requirements; an operator  $\sigma : \mathcal{F} \rightarrow \mathcal{F}$  is a closure operator for the fixpoint-enumeration algorithm iff  $\forall F, G \in \mathcal{F}$ ,  $\sigma$  is

$$\text{extensive: } F \subseteq \sigma(F),$$

$$\text{idempotent: } \sigma(\sigma(F)) = \sigma(F), \text{ and}$$

$$\text{monotonic: } F \subseteq G \Rightarrow \sigma(F) \subseteq \sigma(G).$$

Firstly, extensivity is straightforward to guarantee, we choose  $\sigma(F)$  such that it never removes entities from  $F$ . Secondly, due to idempotency, we require that the closure of a maximal CCPS is the maximal CCPS itself, otherwise it is not a fixpoint and will not be in the output. Thirdly, suppose that the set  $F$  has two supersets that are maximal CCPSSs:  $F', F'' \supseteq F, F' \neq F''$ . Since they are maximal, they both contain an entity that is not in  $F$ , nor in the other maximal CCPS. Extensivity combined with monotonicity forces us to choose  $\sigma(F)$  such that it does not add any entities that are missing from any superset  $G \supseteq F, G \in \mathcal{F}$ .

Hence, we define the closure as follows. Let the set of *compatible entities* be  $\text{Comp}(F) = \{e \in E | F \cup \{e\} \text{ is complete}\}$ , i.e., all entities that can still be added to  $F$ , and let the set of *augmentation entities* be  $\text{Aug}(F) = \{a \in A | F \cup \{a\} \in \mathcal{F}\}$ , i.e., all entities that can be added while leading to a valid CCPS. Then we define the closure operator as in [5]:

$$\sigma(F) = \{e \in \text{Aug}(F) | \text{Comp}(F \cup e) = \text{Comp}(F)\}.$$

This operator is extensive and monotonic, but not idempotent. Without idempotency, we would still enumerate all maximal CCPSSs, but more non-maximal CCPSSs. We achieve idempotency by repeating the closure operator until  $\sigma(F) = F$ . The repetition can be done efficiently by considering only entities that have just become part of  $\text{Aug}(F)$ .

In [5], it is assumed that the dataset does not contain any entity  $e$  that is related to all entities of a neighbouring type, because then all other entities could be in its set of compatible entities ( $\text{Comp}(\{e\}) = E$ ), hence  $\sigma(\emptyset) \supseteq \{e\}$ , while  $e$  need not be part of every CCS. Thus, this assumption is required for the closure operator to be monotonic.

In the current setting, entities that are fully connected to a neighbouring type would not be uncommon and this assumption is not reasonable. For example, there could be a catch-all entity in a hierarchical attribute. Hence, we additionally define  $\sigma(\emptyset) = \emptyset$ . Alternatively, one could redefine  $\text{Comp}$  as  $\text{Comp}(F) = \{e \in E | \exists G \supseteq F \cup \{e\}, G \in \mathcal{F}\}$ , but we leave that to future work. Due to space constraints, we omit the proof that this  $\sigma$  is a closure operator.

**Final Remarks.** The fixpoint-enumeration algorithm enumerates all fixpoints, i.e., any set that results from computing the closure operator. We are only interested in maximal CCPSSs, so we output only those. They are easily identified during the

mining process as they are fixpoints where no entity could be added (Sec. II, Definition 3).

Finally, we allow a user to put any number of constraints on the set of patterns in the form “*any pattern should include at least  $X$  entities of type  $Y$* ”. We implement this by continuously computing upper bounds during the mining process, such that we can prune any branch where the constraints cannot be satisfied any more. A similar approach is followed in [5].

The source code for P-N-RMiner can be found at <https://bitbucket.org/BristolDataScience/p-n-rminer>.

**Computational Complexity.** As stated previously, the number of maximal CCPSs can be exponential in  $|E|$ . Since P-N-RMiner exhaustively enumerates all maximal CCPSs, the worst-case complexity of P-N-RMiner is also exponential in  $|E|$ . Unfortunately, we are not aware of an upper bound on the number of maximal CCPSs, nor do we know the exact worst-case complexity of our algorithm.

It has been shown that the delay time between finding two closed CCSs using the fixpoint-enumeration algorithm is  $O(|E|^3)$  [5]. The algorithm used here is almost the same, except that computing the set of augmentation entities Aug also involves checking the properness constraint. However, the complexity of that is still  $O(|E|)$ , hence the delay time is equivalent. It is not known whether there is a polynomial delay time for maximal CCPSs.

## V. CASE STUDIES

The framework and theory presented in the previous sections give rise to several empirical questions, which we aim to address in this section through three case studies. Our primary contribution is the formalisation of a more general pattern syntax, hence the primary question that we need to verify experimentally is:

- 1) Can we find patterns that are more interesting?

Our secondary contribution is the derivation of an interestingness score that accounts for the dependence between relationship instances of structured attributes. Hence, the second question is:

- 2) Is the new interestingness score relevant?

Thirdly, we present a novel enumeration algorithm. Given the appropriate input, the enumeration algorithm from [6] would output the same maximal CCPSs. However, we claim P-N-RMiner is faster, because it can capitalise on the partial order structure. Hence, the third question is:

- 3) Is the new enumeration algorithm faster?

We aim to answer the first two questions in the following case studies, and also showcase the type of patterns that one can find using the method introduced in this paper. The third question we discuss in Section VI.

**Foursquare Check-Ins.** First we return to the Foursquare Check-ins data discussed in the introduction. This data was gathered by Cheng et al. [11] from several online social media but mostly ( $> 50\%$ ) from Foursquare, and consists of user-ids, check-in times, and venues. The data consists of 225K users and 22M recorded check-ins. Data such as this could be

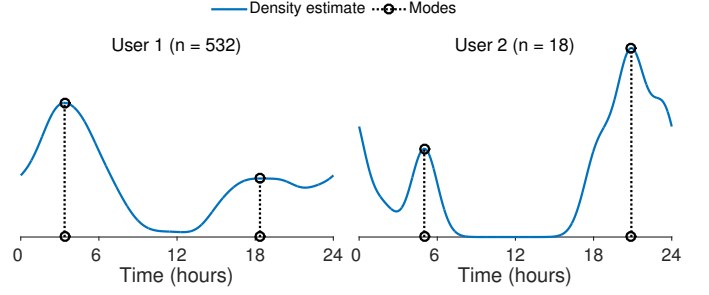


Figure 3. Two examples of density estimation and mode location finding for the check-ins data.

useful to identify patterns of mobility of people, busy times of certain services, etc. Ordinarily, we would represent this data using three entity types and triary relationship instances; user  $x$  checks in at time  $t$  into venue  $y$ . To make this example as simple as possible, we omit the information about venues.

In this case, we are interested in finding patterns in the check-in times across users such as “*many users check in somewhere both between 8.30 and 9.30 in the morning and between 11.30 and 12.30 around noon*”. Such patterns cannot be identified by running P-N-RMiner on the data directly, because relationship instances carry no weights (or any information about their probability). Hence, users that check in frequently and are tracked over a long period of time will have checked in somewhere at many times of the day.

Hence, we preprocess the data by computing kernel density estimates for each user, using a Gaussian kernel with a width of one hour and then locating the modes of their check-in times, with 10-minute precision. Two examples are visualised in Figure 3. As a result, instead of 22M check-ins, the relationship instances correspond to 684K modes, 3 per user on average. This way, more data ensures that our patterns will become more accurate.

We are interested in discovering patterns that possibly include time intervals and not just specific times. As possibilities, we considered asking P-N-RMiner to try intervals up to one, one and a half, and two hours. The reason we consider several options is because the more intervals there are, the more difficult the computational problem is. We identified for each interval size the largest subsampled data that we could run in less than 8 hours<sup>6</sup>, using a reasonable constraint on minimum number of users in any CCPS, each time cutting the data size in half. We found these sample sizes to be  $2^{-8}$  (879 users),  $2^{-9}$  (440 users), and  $2^{-10}$  (220 users), with minimum constraints of 0, 10, and 10 users in all patterns.

Neither of the settings yields substantially more interesting patterns than another. The ‘up to 2-hour intervals’ adds least information to the other two; more than half of the top-100 patterns for that setting contain only intervals that are shorter than 1.5 hours and are thus also present in those results, and the interestingness scores are  $< 0.815$ , while the top-65 for ‘ $\leq 1$ -hour’ and the top-26 for ‘ $\leq 1.5$ -hours’ have higher scores; up to 0.861 and 0.855 respectively. Notice that such scores are not straightforward to interpret, because whether such a score

<sup>6</sup>Unfortunately our current implementation does not use any parallelisation, so it runs only in a single thread.



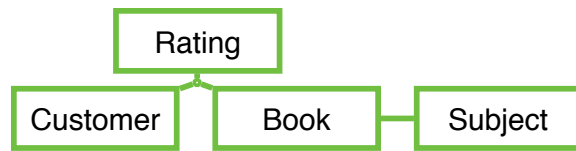


Figure 4. Relational schema of the Amazon Book Ratings data.

is low or high completely depends on the data at hand. For example, the pattern ranked 4<sup>th</sup> for ‘ $\leq 2$ -hours’ is interesting. It contains three intervals and reads: 4.5% of the users checked in frequently between [1.10am–2.30am], [4.30pm–6.30pm], as well as [8.30pm–9.30pm].

The overall most informative pattern that we identify is: 1.6% of the users checked in frequently between [6am–7am], as well as [10.20pm–10.50pm]. This means that, compared to the number of users that check in frequently between those intervals, there is a surprisingly large set of users that checks in frequently during both intervals. This pattern was found in the subsample of 879 users using intervals up to one hour in duration. Interestingly, in that case computing the results without constraints took 2 hours 20 minutes, but all except one pattern in the top-700 (ranked 269) have at least ten users, a result that can be computed in roughly half the time (1h13m).

To confirm that handling intervals is relevant, we identified the top pattern that does not include any intervals; it is ranked 892<sup>nd</sup>, 2962<sup>nd</sup>, and 10138<sup>th</sup>, for the three cases respectively. This shows beyond any doubt that patterns with intervals are more interesting. We also test the relevance of the new interestingness score, by comparing the ranking of P-N-RMiner against N-RMiner on data augmented such that they produce the same patterns. We find that Kendall’s tau is 0.337 and 0.352, respectively (N-RMiner did not finish in time on the third dataset), which highlights that accounting for the partial order when computing interestingness is highly relevant.

**Amazon Book Ratings.** As a second case study, we downloaded a snapshot of Amazon product reviews from SNAP<sup>7</sup>. This dataset contains around 500K products, 8M reviews with ratings from 1 to 5, and 2.5M product category memberships. From this we selected all reviews about books and uniformly sub-sampled 1% of the customers.

Every book has multiple category memberships which are given as paths in the Amazon product category hierarchy. From this we extracted the relationship between books and categories and the hierarchy itself, keeping two levels below the category *Book*  $\rightarrow$  *Subject*. The dataset that we obtain has the structure shown in Figure 4 and consists of 22,003 books, 9,855 customers, 417 hierarchically structured book subjects, as well as 36,415 ratings and 53,403 subject memberships.

We ran P-N-RMiner on this dataset with constraints of at least 6 books and 20 customers. As an example, we present the most highly-ranked pattern. This contains 23 customers and 8 books, all of which are different versions of the book “Left Behind: A Novel of the Earth’s Last Days”, a rating [1–5] and the subjects *Fiction* and *Christianity*. To our surprise, we found that most of the patterns in the result are like this; different versions of the same book (hard cover, audiobook, etc.).

Inspection of the raw data led us to the hypothesis that this happens because reviews are copied across different versions of the same book. Unfortunately, the text of reviews was not crawled, so it is not straightforward to identify reviews for different items that are equivalent. We attempted to tackle this problem by keeping only one such version of a book by looking for reviews that have the same date, rating, and user. However, after removing duplicates using this procedure, it appears that little structure remains in the data.

We also ran N-RMiner on the same dataset, augmenting it with all the implied relationship instances. We see that the same pattern is now ranked at the 21<sup>st</sup> position. This is because N-RMiner does not take into account the dependencies between the intervals and, as a result, intervals are by definition more highly connected and relationship instances containing intervals are more probable. This confirms that our new derivation of the interestingness score is indeed relevant.

**Fisher’s Iris Data.** The Iris data<sup>8</sup> has been pervasively used in machine learning and pattern recognition text books. The data consists of 150 measurements of plants. Each has four numerical attributes and a class label (one of three species). In Section II, we have shown that P-N-RMiner can be used to mine tiles and frequent patterns. However, it can also be used to mine subgroups and subspace clusters, which we highlight in this case study.

*Subgroup discovery* is a form of pattern mining where a user chooses a target attribute and the aim is to find rules that predict high values of this attribute (or rules that predict *true* if the attribute is binary). For the Iris data, this means that we would like to find rules based on the four numerical attributes that predict a specific class label. We model the data as five entity types. We discretise each numerical attribute to ten different values using equal spacing and include intervals up to six adjacent values. This substantially reduces the computation time, while hardly affecting the patterns.

We then ran P-N-RMiner with a constraint that all patterns have to include a class label. The top pattern for each class is visualised in Figure 5. Interestingly, all top patterns include values for all four numerical attributes, indicating that they are all informative for the class label and the set of points that they describe. The first pattern that omits an attribute is ranked 120<sup>th</sup> and is equivalent to the second most informative pattern in the data (and second most informative for class 1), except that it omits *sepal width*. Figure 5 visually confirms that *sepal width* is the least informative feature for that pattern.

*Subspace clustering* is a form of pattern mining that is unsupervised. The goal is to discover clusters in the data, but unlike traditional clustering, the goal is not to provide a full partitioning of the data, and there is no requirement to use all variables. Our framework has roughly the same aim, and could as such be considered a relational (exhaustive) approach to subspace clustering. Like in the case of the check-ins data, our framework enables identification of patterns that are otherwise unattainable using existing methods.

To find subspace clusters in the data, we ran P-N-RMiner without constraints on the Iris data without the class labels. As output we find 25,365 patterns. The top pattern

<sup>7</sup><https://snap.stanford.edu/data/amazon-meta.html>

<sup>8</sup><https://archive.ics.uci.edu/ml/datasets/Iris>



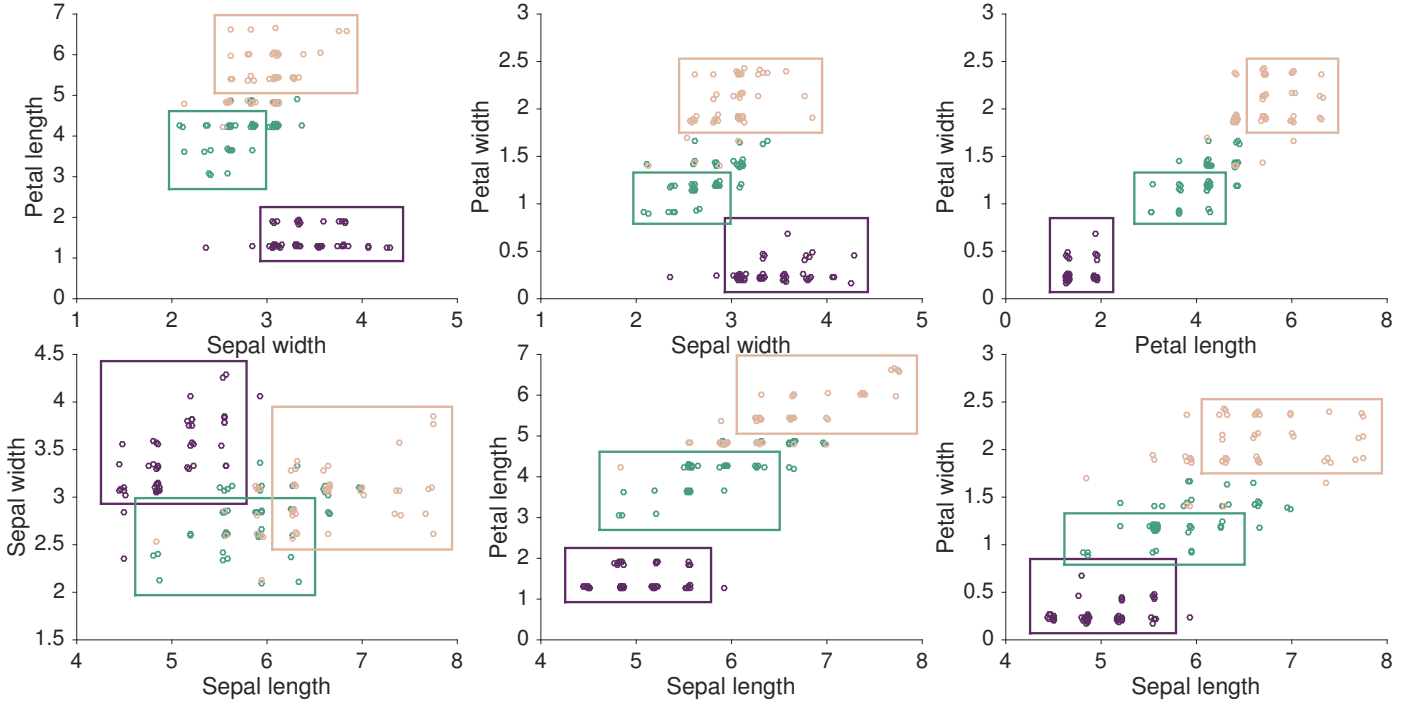


Figure 5. Visualisation of the full Iris data, projected for each pair of features. Colors depict class labels and the boxes represent the top subgroup pattern for each class, as discovered by P-N-RMiner. Incidentally, each most informative pattern includes all four attributes.

is:  $pl = [1.295-1.885]$ ,  $pw = [0.22-0.7]$ ,  $sl = [4.84-5.2]$ ,  $sw = [3.08-4.04]$ , with an interestingness score of 1.4744. So, it is similar to the top pattern predicting class 1 (see Figure 5), except that the intervals for *sepal length* and *sepal width* are slightly more narrow. The first *subspace* cluster occurs at rank 347 and is quite specific already:  $pl = [1.295-1.885]$ ,  $pw = [0.22-0.7]$ ,  $sl = [4.48-5.2]$ , with an interestingness score of 1.1040, again omitting *sepal width*.

As a final remark, we are not suggesting that P-N-RMiner can replace all existing subgroup discovery and subspace clustering methods, because P-N-RMiner has high computational cost owing to the exhaustive search strategy.

## VI. SCALABILITY

To test the scalability of the algorithm and study what we gain by using the attribute structure in the form of the properness constraint, we again look at the Foursquare check-ins data. We created 11 versions of the data, each time throwing away half of the users and their relationship instances (the check-in modes). We then ran both N-RMiner [6] on augmented data with the additional entities and relationship instances, and P-N-RMiner, which then find exactly the same set of CCPSSs.

We are interested also in how the depth of the partial order of an attribute affects the scalability and potential speed-up by P-N-RMiner. Hence, we tested runtimes for 6, 9, and 12 levels, i.e., time intervals up to one, one and a half, and two hours. We exhaustively tested constraints on the number of users from 10, 20, 40, etc. up to the sample size. We stopped any experiment that had not finished after 24 hours.

A comparison of runtimes for all cases where N-RMiner finished successfully is given in Figure 6. The general trend

is that P-N-RMiner is faster in 71 out of this subset of 100 experiments, and that the speed-up grows as the computation time grows. The largest observed speed-up is a factor of 8; 12.3 vs. 1.5 hours, for the full data, 12 levels, and mining patterns with at least 81920 users. N-RMiner is mainly faster (up to a factor of 2) for runtimes shorter than 1 second. Not shown in the figure is that P-N-RMiner uses substantially less memory. For example, for the full data with 12 levels, P-N-RMiner uses 5.5 MB of memory at its peak, while N-RMiner uses more than 10 GB.

Runtimes of P-N-RMiner for increasing sample sizes are illustrated in Figure 7. It appears that the number of levels, i.e., the depth of the partial order, actually only has a small effect on the runtime—for example, compare the trend of brown points across the depths. Yet, the number of patterns in the data explodes much more easily. For depth 12, we can only enumerate all patterns in the samples that have at most 879 users (0.4% of the data). There appears to be a linear relationship between the size of the data and the runtime if the number of patterns is equivalent; for any of the subfigures, one could fit a straight line through measurements that have roughly the same number of patterns (i.e., points of the same color). This relationship holds even when there are no patterns.

## VII. RELATED WORK

**Exploratory vs. Predictive Patterns.** The broad purpose of the framework presented in this paper is to facilitate *exploration* of data in an entirely unsupervised manner. This distinguishes the framework from other types of local pattern for multi-relational data mining such as Safarii [4], and more generally from approaches based on inductive logic programming. These alternative frameworks operate by a user selecting

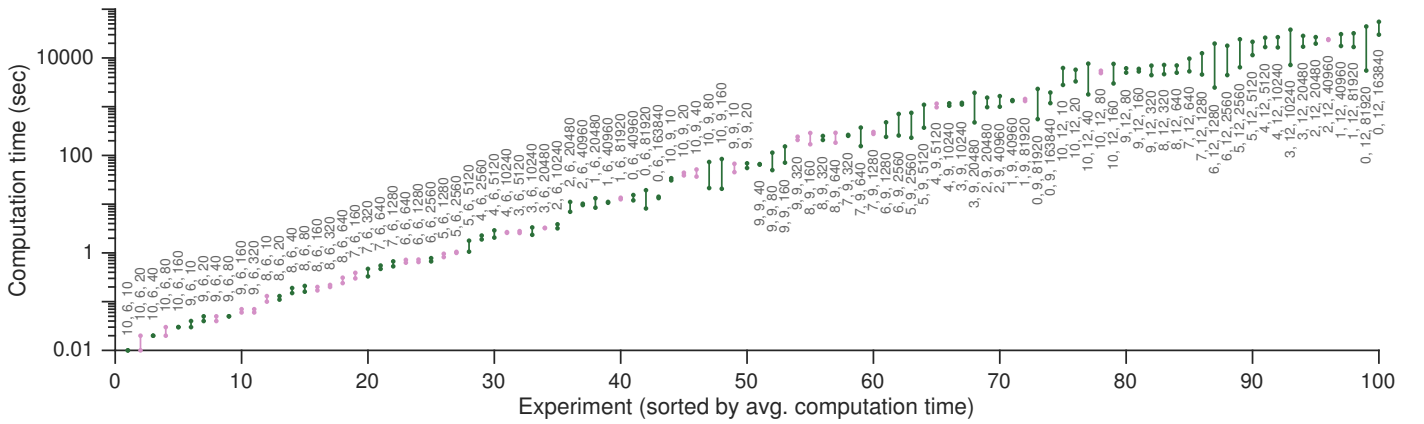


Figure 6. Comparison of computation times between P-N-RMiner (introduced here) and NR-Miner [6]. On the x-axis are experiments, sorted by average computation time over the two methods, and on the y-axis is computation time. The runtime of each experiment is visualised by two dots, corresponding to the runtime of the two methods. For dark/green bars and dots, P-N-RMiner is the lower dot of the two (and thus faster), while for light/pink bars and dots, N-RMiner is lower dot of the two (and thus faster). P-N-RMiner is typically faster, but especially for problems that are computationally more demanding.

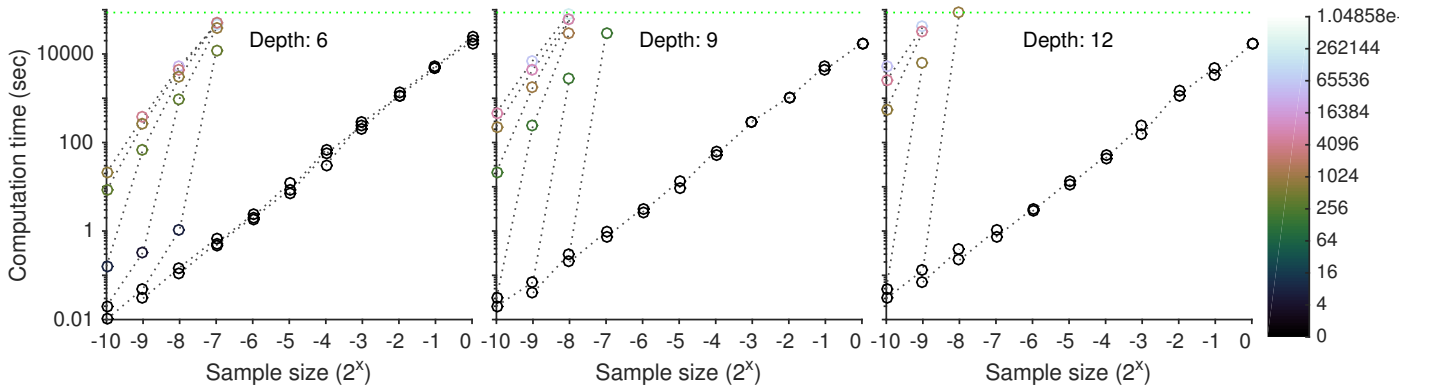


Figure 7. Scalability of P-N-RMiner for increasing data size. Each dotted line with circles corresponds to a constraint on the minimum number of users. For any such constraint there are typically two or three measurements (sample sizes) that finished in 24 hours, due to the constraint being too low and demanding, or too high (larger than the sample size). On the x-axis is the sample size, on the y-axis the runtime. The color of a circle corresponds to the number of patterns (CCPSs) enumerated, ranging from 0 (black) to ca. 250k (light blue), see also the colorbar. There appears to be linear relation between runtime and sample size when the number of patterns is the same, and the depth of the partial order structure has a limited effect as well.

one or a set of attributes as a target, after which an algorithm builds rules to predict that target using the full relational data. Wu et al. [12] introduced a method for finding interesting chains of biclusters in relational data, which has a similar goal as our framework. Their approach differs in that they only consider binary relationships, they employ a heuristic greedy algorithm to find interesting patterns, and their method does not account for structure of attributes in any way.

**Pattern Syntax.** The pattern syntax proposed in this paper is unique in being both relational and able to deal with structured attribute types such as ordinal and real-valued attributes, taxonomy terms, and more. The proposed pattern syntax, in being *local*, owes to the frequent pattern mining literature. Indeed, the CCS pattern syntax [5], which it generalises, has already been shown to be a generalisation itself of a local pattern type in binary databases known as *tiles* [13], which are essentially equivalent to frequent itemsets.

**Structured Attribute Types.** Real-valued and ordinal attributes have also been dealt with before in local pattern mining, in subgroup discovery and exceptional model mining. For example, in subgroup discovery, approaches have been de-

veloped to infer subgroup descriptions in terms of intervals for real-valued attribute types and subsets of categorical attributes. A notable paper in this regard is [14], where an efficient algorithm is introduced for finding optimal subgroups using any convex quality measure. Exceptional model mining, on the other hand, aims to extend subgroup discovery beyond a single target attribute [15]. None of these approaches, however, are as generic as our proposed approach: they are either ad hoc or remain limited to a very specific types of structured attributes. The approach of modelling the structure of the attributes as a partial order is also entirely novel.

**Interestingness Formalisations.** The formalisation of interestingness of local patterns is a highly active research area, with most research targeted on itemsets in binary databases. This makes sense, as the problem is most acute for exploratory data mining approaches, in the absence of a particular set of target attributes to be predicted. Many approaches to formalising interestingness are based on modelling the *unexpectedness* of a pattern: the extent to which the pattern presents novel, surprising, or unexpected information to the user. A recent survey is [16].

There are three major lines of research aimed at mining (sets of) interesting local patterns. Constrained randomisation techniques are based on the assumption that a pattern is more interesting if it is not present in randomised data [17], [18], [19]. Methods based on the Minimum Description Length principle assume that a pattern is more interesting if it provides better compression [20]. Approaches based on the Maximum Entropy (MaxEnt) principle assume a pattern is more interesting the more surprising it is given a MaxEnt-based background model [8], [9]. Both randomisation and MaxEnt approaches have been shown to allow for accounting prior knowledge, enabling subjective interestingness and iterative data mining.

The MaxEnt approach has been shown to be highly flexible in terms of pattern types [21]. Additionally, it has been used successfully to quantify interestingness patterns for the framework that we directly build upon [5]. For these reasons we used this paradigm to formalise the interestingness of the patterns in the current paper. Clearly, a direct application of interestingness as defined in [5] would not have yielded desirable results, as the dependencies between relationship instances would be ignored (see also the empirical results).

**Enumeration Algorithms.** The algorithm derived for enumerating all maximal CCPs is based on the generic fixpoint-enumeration algorithm for enumerating all closed sets in a strongly accessible set system, introduced by Boley et al. [10]. This algorithmic scheme has been used before in the data mining literature for enumerating maximal CCSs [5], including extensions to  $n$ -ary relations [6] and approximate CCSs [7]. Here, we adhere to the same algorithmic scheme. In order to be able to use the scheme, we model the structure of attributes as a partial order, augment the pattern syntax, and add a properness constraint to the definition of the set of augmentation elements. As may be apparent, these changes are not trivial, and neither is the proof that the algorithmic scheme still works.

## VIII. CONCLUSIONS

An important obstacle for the adoption of exploratory data mining techniques in general, and local pattern mining approaches in particular, is their limited flexibility in terms of data type to which they can be applied (e.g., only tabular data), and type of pattern they can generate, e.g. subgroups, itemsets,  $n$ -sets. In reality, however, data is often complexly structured (as in, e.g., a relational database), and additionally there is often structure among the different values data attributes may attain, i.e., attribute values can be ordinal, interval, taxonomy terms, and more.

Attempts to resolve this inflexibility for specific data and pattern types are numerous. Yet, we are unaware of any generic approach that comes close to subsuming the range of pattern syntaxes considered by the local pattern mining research community, allowing for data types of a broad range of structures. The contributions in the present paper may be an important step in this direction.

Our contributions raise a number of new research challenges. Ideally, the pattern syntax is tolerant to missing relations to ensure noise resilience, similar to [7]. The interestingness can be made more versatile by considering a more varied range of prior belief types. Another interesting question is whether the enumeration algorithm could still be improved.

Our algorithm is similar to the Bron-Kerbosch algorithm for enumerating maximal cliques in a graph, for which it is known that the worst case complexity of  $O(3^{n/3})$  is optimal, since it is equivalent to the number of maximal cliques in a graph [22]. Yet another interesting direction for future work is developing heuristic algorithms for finding interesting CCPs directly, in order to avoid the costly exhaustive search step.

*This work was supported by the European Union (ERC Grant FORSID 615517) and the EPSRC (Grant EP/M000060/1).*

## REFERENCES

- [1] C. Borgelt, "Frequent item set mining," *WIREs: DMKD*, vol. 2, no. 6, pp. 437–456, 2012.
- [2] L. Cerf, J. Besson, C. Robardet, and J.-F. Boulicaut, "Data peeler: Constraint-based closed pattern mining in  $n$ -ary relations," in *Proc. of SDM*, 2008, pp. 37–48.
- [3] F. Herrera, C. J. Carmona, P. González, and M. J. del Jesus, "An overview on subgroup discovery: foundations and applications," *KAIS*, vol. 29, no. 3, pp. 495–525, 2011.
- [4] A. J. Knobbe, *Multi-relational data mining*. IOS Press, 2006.
- [5] E. Spyropoulou, T. De Bie, and M. Boley, "Interesting pattern mining in multi-relational data," *DMKD*, vol. 28, no. 3, pp. 808–849, 2014.
- [6] —, "Mining interesting patterns in multi-relational data with  $N$ -ary relationships," in *Proc. of DS*, 2013, pp. 217–232.
- [7] E. Spyropoulou and T. De Bie, "Approximate multi-relational patterns," in *Proc. of DSAA*, 2014, pp. 477–483.
- [8] T. De Bie, "An information-theoretic framework for data mining," in *Proc. of KDD*, 2011, pp. 564–572.
- [9] —, "Maximum entropy models and subjective interestingness: an application to tiles in binary databases," *DMKD*, vol. 23, no. 3, pp. 407–446, 2011.
- [10] M. Boley, T. Horváth, A. Pogné, and S. Wrobel, "Listing closed sets of strongly accessible set systems with applications to data mining," *TCS*, vol. 411, no. 3, pp. 691–700, 2010.
- [11] Z. Cheng, J. Caverlee, K. Lee, and D. Z. Sui, "Exploring millions of footprints in location sharing services," in *Proc. of ICWSM*, 2011, pp. 81–88.
- [12] H. Wu, J. Vreeken, N. Tatti, and N. Ramakrishnan, "Uncovering the plot: detecting surprising coalitions of entities in multi-relational schemas," *DMKD*, vol. 28, no. 5–6, pp. 1398–1428, 2014.
- [13] F. Geerts, B. Goethals, and T. Mielikäinen, "Tiling databases," in *Proc. of DS*, 2004, pp. 278–289.
- [14] M. Mampaey, S. Nijssen, A. Feelders, and A. Knobbe, "Efficient algorithms for finding richer subgroup descriptions in numeric and nominal data," in *Proc. of ICDM*, 2012, pp. 499–508.
- [15] D. Leman, A. Feelders, and A. Knobbe, "Exceptional model mining," in *Proc. of ECML-PKDD*, 2008, pp. 1–16.
- [16] K.-N. Kontonassios, E. Spyropoulou, and T. De Bie, "Knowledge discovery interestingness measures based on unexpectedness," *WIREs DMKD*, vol. 2, no. 5, pp. 386–399, 2012.
- [17] A. Gionis, H. Mannila, T. Mielikäinen, and P. Tsaparas, "Assessing data mining results via swap randomization," *TKDD*, vol. 1, no. 3, p. 14, 2007.
- [18] M. Ojala, N. Vuokko, A. Kallio, N. Haiminen, and H. Mannila, "Randomization of real-valued matrices for assessing the significance of data mining results," in *Proc. of SDM*, vol. 8, 2008, pp. 494–505.
- [19] J. Lijffijt, P. Papapetrou, and K. Puolamäki, "A statistical significance testing approach to mining the most informative set of patterns," *DMKD*, vol. 28, no. 1, pp. 238–263, 2014.
- [20] J. Vreeken, M. van Leeuwen, and A. Siebes, "Krimp: mining itemsets that compress," *DMKD*, vol. 23, no. 1, pp. 169–214, 2011.
- [21] T. De Bie, "Subjective interestingness in exploratory data mining," in *Proc. of IDA*, 2013, pp. 19–31.
- [22] E. Tomita, A. Tanaka, and H. Takahashi, "The worst-case time complexity for generating all maximal cliques and computational experiments," *TCS*, vol. 363, no. 1, pp. 28–42, 2006.